

Fundamental Limits of Coded Caching: Improved Delivery Rate-Cache Capacity Trade-off

Mohammad Mohammadi Amiri, *Student Member, IEEE* and Deniz Gündüz, *Senior Member, IEEE*

Abstract—A *centralized coded caching* system, consisting of a server delivering N popular files, each of size F bits, to K users through an error-free shared link, is considered. It is assumed that each user is equipped with a local cache memory with capacity MF bits, and contents can be proactively cached into these caches over a low traffic period; however, without the knowledge of the user demands. During the peak traffic period each user requests a single file from the server. The goal is to minimize the number of bits delivered by the server over the shared link, known as the *delivery rate*, over all user demand combinations. A novel coded caching scheme for the cache capacity of $M = (N - 1)/K$ is proposed. It is shown that the proposed scheme achieves a smaller delivery rate than the existing coded caching schemes in the literature when $K > N \geq 3$. Furthermore, we argue that the delivery rate of the proposed scheme is within a constant multiplicative factor of 2 of the optimal delivery rate for cache capacities $1/K \leq M \leq (N - 1)/K$, when $K > N \geq 3$.

Index Terms—Centralized coded caching, network coding, proactive caching.

I. INTRODUCTION

The growing number of users and their increasing appetite for high data rate content leads to network traffic congestion, particularly during peak traffic periods, whereas the resources are often underutilized during off-peak periods. Exploiting increasingly low-cost and abundant local storage capacity to proactively cache content at user devices is an effective way to utilize channel resources during off-peak hours, and mitigate the burden of high network load at times of heavy demand [1], [2].

In order to model this dichotomy between peak and off-peak traffic periods, recent works on coded content caching consider two phases: In the *placement phase*, which corresponds to periods of low network traffic, the cache memory of each user is filled by a central server without the knowledge of users' future demands. The main limitation of this phase is the capacity of users' caches. All user requests are revealed simultaneously during the peak traffic period. It is assumed that each user requests a single file from among a finite database of popular contents. Then the *delivery phase* follows, in which a common message is transmitted to all the users in the system over an error-free shared channel. Each user tries to reconstruct the file it requests using its local cache content

as well as the bits delivered by the server over the shared link. The goal of the server is to make sure that all the user demands, no matter which files are requested by the users, are satisfied at the end of the delivery phase. For a given number of files in the database, and given cache sizes at the users, the minimum rate, referred to as *delivery rate*, at which data must be delivered through the shared link, independently of users' particular demands, is considered as the performance measure of a coded caching algorithm. Our goal is to design the placement and delivery phases jointly in order to minimize the delivery rate.

In the case of uncoded caching, parts of popular contents are stored in the local cache memories, and once the user requests are revealed, remaining parts are delivered by the server over the shared link. The corresponding gain relative to not having a cache is called *local caching gain*, and depends on the local cache capacity [3], [4]. On the other hand, it has been shown in [5] that *coded caching* provides a *global caching gain*, where users benefit not only from their own local cache, but also from the available cache memory across the network. Coded caching provides a novel method to mitigate network congestion during peak traffic hours by creating and exploiting coded multicasting opportunities across users.

In a *centralized coded caching* scheme, it is assumed that the central server knows the exact number of users in the system, and carefully places contents in the user caches during off-peak hours. A novel centralized coded caching scheme for a network of K users requesting N popular files of the same size is proposed in [5], which has been shown to be optimal when the cache placement phase is uncoded and $K \geq N$ [6]. Authors in [7] consider an alternative coded caching scheme, which was originally proposed in [5] for three users, and show that it is optimal when the number of users, K , is not less than the number of popular files, i.e., $N \leq K$, and the normalized cache capacity M satisfies $M \leq 1/K$, i.e., a relatively small cache size. The delivery rate is further improved by the coded caching schemes investigated in [8] and [9]. Theoretical lower bounds on the delivery rate have also been derived to characterise the optimal performance of a caching system [5], [10]–[12]. In general, the minimum delivery rate for coded caching remains an open problem even in the symmetric setting considered in the aforementioned previous works.

The scheme of [5] has been extended to the decentralized setting in [13], in which the identity of the users requesting files during the delivery phase are not known in advance to perform the placement in coordination across caches. It has been further extended to multi-layer caching [14], caching files with distinct sizes [15] and popularities [16], [17], caching

The authors are with Imperial College London, London SW7 2AZ, U.K. (e-mail: m.mohammadi-amiri15@imperial.ac.uk; d.gunduz@imperial.ac.uk).

This work is funded partially by the European Research Council Starting Grant project BEACON (Project number 677854), and by the European Union's Horizon 2020 Research and Innovation Programme under grant agreement 690893, project TACTILENet: Towards Agile, effiCient, auTonomous and massIvely LargeE Network of things.

to users with distinct cache capacities [18], [19], hierarchical coded caching [14], [20], caching files with lossy reconstruction [21]–[23], online cache placement [24], and delivery over a noisy shared link [25]–[27]. Similar caching techniques have also been employed in various other applications, e.g., device-to-device caching [28], [29] and femtocaching [30], [31].

In this paper, we build upon our previous work in [9], and propose a novel centralized coded caching scheme, when the cache capacity of the users is given by $M = (N - 1)/K$. This new caching scheme utilizes coded content placement, in which contents are partitioned into smaller chunks, and pairwise XOR-ed contents are placed in the user caches. The delivery phase utilizes both coded and uncoded transmission. We show that the proposed caching scheme requires a smaller delivery rate (evaluated for the worst-case user demands) compared to the best achievable scheme in the literature for the same cache capacity, when $N < K$. We then extend the improvement in the delivery rate to a larger range of cache capacities utilizing the memory-sharing argument. Finally, we show that the delivery rate achieved by the proposed caching scheme is within a constant multiplicative factor of 2 of the optimal delivery rate for cache capacities satisfying $1/K \leq M \leq (N - 1)/K$, when $K > N \geq 3$. We believe that the ideas behind our centralized coded caching scheme may lead to further improvements on the delivery rate in decentralized as well as online caching systems.

We remark that the proposed caching scheme improves the performance upon the state-of-the-art when there are more users in the system than the number of files in the database, i.e., $N < K$, and cache capacity of each user is relatively small. This scenario is valid for contents that become highly popular over the Internet, and are demanded by a huge number of users, each equipped with a cache memory of comparatively small size, in a relatively short time interval, for example, viral videos distributed over social networks, new episodes of popular TV series, breaking news videos, or for broadcasting different software updates to millions of clients.

The rest of this paper is organized as follows. In Section II, we investigate the system model, and present the relevant previous results in the literature. In Section III, a novel centralized coded caching scheme is proposed, and its delivery rate is analyzed and compared with the state-of-the-art results both theoretically and numerically. All the proofs can be found in the Appendix. Finally, conclusions are included in Section IV.

II. SYSTEM MODEL AND PREVIOUS RESULTS

We assume that there is a central server broadcasting data to $K \in \mathbb{Z}^+$ users, U_1, \dots, U_K , through a shared error-free link, where \mathbb{Z}^+ is the set of positive integers. As depicted in Fig. 1, the server has $N \in \mathbb{Z}^+$ files in its database, each with the size of F bits, that is, file n , denoted by W_n , for $n = 1, \dots, N$, is a random variable uniformly distributed over $[2^F] \triangleq \{1, \dots, 2^F\}$. We denote the whole database by $\mathbf{W} \triangleq (W_1, \dots, W_N)$. Each user is assumed to have a local cache of capacity MF bits.

Similarly to [5], the system operates in two phases. In the *placement phase*, each user's cache is filled with bits

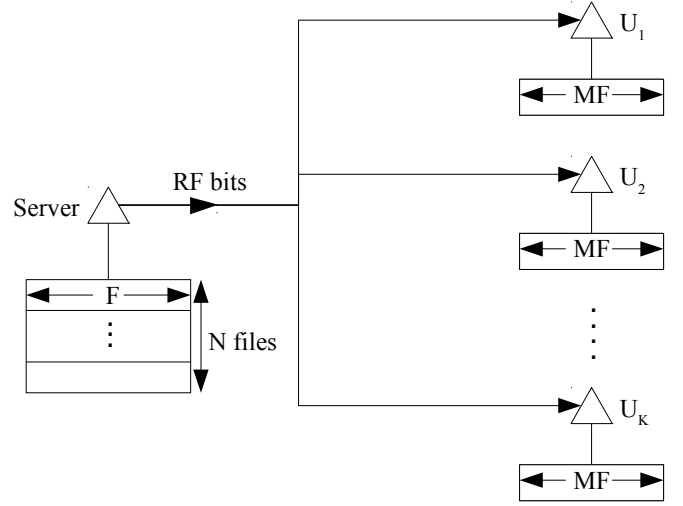


Fig. 1. Illustration of a centralized caching system consisting of a server with a database of N popular files, each with size F bits, serving K users, each with a cache of capacity MF bits, requesting a single file from the database. These requests are served simultaneously through an error-free shared link.

depending on the whole database and the capacity of the user caches, M . The content of user k 's cache at the end of the placement phase is denoted by Z_k . At the end of the placement phase, each user requests one of the files from the database. We use $\mathbf{d} = (d_1, \dots, d_K)$ to denote the demand vector, where $d_k \in [N]$ corresponds to the demand of user U_k . In the *delivery phase*, the server, having received the requests of all the users, transmits a common message of size RF bits over the shared link. Note that, in the centralized caching model considered here, this common message depends not only on the contents in the database and the user requests, but also on the contents of users' caches. At the end of the delivery phase, each user k tries to decode its requested file W_{d_k} using its cache content Z_k , and the common message received over the shared link.

Definition. An (M, R, F) caching and delivery code for the above caching system with K users and N files consists of

- i) K caching functions:

$$\phi_k : [2^F]^N \rightarrow [2^{\lfloor FM \rfloor}], \quad (1)$$

which maps the database \mathbf{W} to the cache content Z_k of user k , i.e., $Z_k = \phi_k(\mathbf{W})$;

- ii) delivery encoding function:

$$\psi : [2^F]^N \times [N]^K \rightarrow [2^{\lfloor FR \rfloor}], \quad (2)$$

which maps the database \mathbf{W} and the particular demand vector \mathbf{d} to a message X over the shared link, i.e., $X = \psi(\mathbf{W}, \mathbf{d})$;

- iii) K decoding functions:

$$\mu_k : [2^{\lfloor FM \rfloor}] \times [2^{\lfloor FR \rfloor}] \times [N]^K \rightarrow [2^F], \quad (3)$$

which maps the cache content Z_k , message over the shared link X , and the demand vector \mathbf{d} to the reconstructed message \hat{W}_k at user k , i.e., we have $\hat{W}_k = \mu_k(Z_k, X, \mathbf{d})$.

The probability of error of an (M, R, F) caching and delivery code is defined as

$$P_e \triangleq \max_{(d_1, d_2, \dots, d_K)} \Pr \left\{ \bigcup_{k=1}^K \{ \hat{W}_k \neq W_{d_k} \} \right\}. \quad (4)$$

In this model, M is the normalized cache capacity while R is the delivery rate, which corresponds to the number of bits transmitted over the shared link, also normalized by the file length F .

Definition. The delivery rate-cache capacity pair (R, M) is *achievable* if for any $\varepsilon > 0$, there exist a large enough F and a corresponding (M, R, F) caching and delivery code with $P_e < \varepsilon$.

There is a trade-off between the cache capacity of the users and the delivery rate. For example, when $M = 0$, in the worst case, users request as distinct files as possible, and the server has to transmit all the requests over the shared link, that is, the delivery rate has to be at least $R = \min \{N, K\}$. In the other extreme case, when the cache capacities are large enough to store all the N files, i.e., when $M = N$, all the requests can be satisfied directly from the local caches, and the delivery rate can be zero, i.e., $R = 0$. In general, we define the delivery rate-cache capacity trade-off $R^*(M)$ as follows:

$$R^*(M) \triangleq \inf \{R : (R, M) \text{ is achievable}\}. \quad (5)$$

Our goal is to obtain the delivery rate-cache capacity trade-off for all possible cache capacity values in between the above extreme scenarios by designing the placement and delivery phases jointly.

When $N > K$, the centralized coded caching scheme proposed by Maddah-Ali and Niesen in [5] for cache capacities $M = tN/K$, for $t = 1, \dots, K$, is the best known achievable scheme in the literature. On the other hand, when $N \leq K$, the scheme proposed in [7] is shown to be optimal for $M \leq 1/K$. In this case, to characterize the best achievable delivery rate in the literature for $M \geq 1/K$, we calculate the delivery rate achieved by memory-sharing between the schemes proposed for $M = 1/K$ in [7] and for $M = tN/K$ in [5], for $t \in [K]$, as follows:

$$R_{M,t} = \alpha_t M + \beta_t, \quad \text{for } \frac{1}{K} \leq M \leq \frac{tN}{K}, \quad (6)$$

where

$$\alpha_t = \frac{K(K-t)}{(t+1)(tN-1)} - \frac{N(K-1)}{tN-1}, \quad (7a)$$

$$\beta_t = N - \frac{N}{K} - \frac{K-t}{(t+1)(tN-1)} + \frac{N(K-1)}{K(tN-1)}. \quad (7b)$$

The value of $t \in [K]$ that minimizes α_t in (7a) is denoted by t^* :

$$t^* \triangleq \arg \min_{t \in [K]} \{\alpha_t\}. \quad (8)$$

This leads to the straight line R_{M,t^*} in (6) with the lowest slope, which characterizes a range of delivery rates achieved through memory-sharing between the schemes proposed for $M = 1/K$ in [7] and for $M = t^*N/K$ in [5]. This scheme, referred to as the Maddah-Ali-Niesen-Chen (MNC) scheme, is considered as the state-of-the-art for $1/K \leq M \leq t^*N/K$

	Cache content
User 1	$W_{1,1} \oplus W_{2,1}, W_{2,1} \oplus W_{3,1}$
User 2	$W_{1,2} \oplus W_{2,2}, W_{2,2} \oplus W_{3,2}$
User 3	$W_{1,3} \oplus W_{2,3}, W_{2,3} \oplus W_{3,3}$
User 4	$W_{1,4} \oplus W_{2,4}, W_{2,4} \oplus W_{3,4}$
User 5	$W_{1,5} \oplus W_{2,5}, W_{2,5} \oplus W_{3,5}$

Fig. 2. Cache placement for the proposed coded caching scheme for $N = 3$ files and $K = 5$ users, each equipped with a cache of capacity $M = 2/5$.

throughout this paper, and achieves a delivery rate of

$$R_b(M) = R_{M,t^*} = \alpha_{t^*} M + \beta_{t^*}, \quad \text{for } \frac{1}{K} \leq M \leq \frac{t^*N}{K}. \quad (9)$$

For $M \geq t^*N/K$, Maddah-Ali and Niesen's scheme in [5] again achieves the best performance in the literature.

In the following, we will introduce the placement and delivery phases for the proposed caching strategy for a cache capacity of $M = (N-1)/K$. We will show that, for this particular cache capacity the proposed scheme achieves a smaller delivery rate than the state-of-the-art presented above. We then characterize a delivery rate-cache capacity trade-off using memory-sharing between our scheme and the existing ones, which extends the improvement to a larger range of cache capacity values.

III. PROPOSED CODED CACHING SCHEME

Before we present a detailed description and analysis of the proposed coded caching scheme, we illustrate it on a simple example highlighting its main ingredients. This example will allow us not only to provide the intuition behind the proposed caching scheme, but also to show its superiority over the MNC scheme.

Example 1. Consider a caching system with a database of $N = 3$ files, W_1 , W_2 and W_3 . There are $K = 5$ users in the system, each of which is equipped with a cache of capacity $M = (N-1)/K = 2/5$. To perform the placement phase, each file W_i , $\forall i$, is first divided into $K = 5$ non-overlapping subfiles $W_{i,j}$, each of the same length $F/5$ bits, for $j = 1, \dots, 5$. The following contents are then cached by user U_k , $\forall k$, in the placement phase:

$$Z_k = (W_{1,k} \oplus W_{2,k}, W_{2,k} \oplus W_{3,k}), \quad (10)$$

where \oplus is the bitwise XOR operation. Since each subfile $W_{i,j}$ has a length of $F/5$ bits, the cache placement phase satisfies the memory constraint. See Fig. 2 for an explicit illustration of the cache contents at the end of the placement phase.

We argue that for the proposed caching scheme with $N < K$, the worst-case user demands happens when each file is requested by at least one user. This fact will later be clarified in Remark 1. By re-labeling the files and re-ordering the users,

without loss of generality, the user demand combination is assumed to be $\mathbf{d} = (1, 1, 1, 2, 3)$.

In the delivery phase, each subfile $W_{i,j}$, $\forall i, j$, is further divided into $N - 1 = 2$ distinct pieces $W_{i,j}^{(l)}$, for $l = 1, 2$, each of size $F/10$ bits, i.e., $W_{i,j} = (W_{i,j}^{(1)}, W_{i,j}^{(2)})$, $\forall i, j$. Accordingly, cache contents (10) can be rewritten as

$$Z_k = \bigcup_{l=1}^2 (W_{1,k}^{(l)} \oplus W_{2,k}^{(l)}, W_{2,k}^{(l)} \oplus W_{3,k}^{(l)}). \quad (11)$$

The contents are then delivered by the server in three different parts. The following contents are sent in each part of the delivery phase:

Part 1: $W_{2,1}^{(1)}, W_{3,1}^{(2)}, W_{2,2}^{(1)}, W_{3,2}^{(2)}, W_{2,3}^{(1)}, W_{3,3}^{(2)}, W_{1,4}^{(1)}, W_{3,4}^{(2)}, W_{1,5}^{(1)}, W_{2,5}^{(2)}$,

Part 2: $W_{1,1} \oplus W_{1,2}, W_{1,2} \oplus W_{1,3},$

Part 3: $W_{2,1}^{(2)} \oplus W_{2,2}^{(2)}, W_{2,2}^{(2)} \oplus W_{2,3}^{(2)}, W_{1,4} \oplus W_{2,3}^{(2)}, W_{3,1}^{(1)} \oplus W_{3,2}^{(1)}, W_{3,2}^{(1)} \oplus W_{3,3}^{(1)}, W_{1,5}^{(2)} \oplus W_{3,3}^{(1)}, W_{2,5}^{(1)} \oplus W_{3,4}^{(1)}.$

Having received the contents delivered in part 1, each user can retrieve all the subfiles placed in its own cache in XOR-ed form. For example, user U_1 can decode all the subfiles $W_{i,1}$, for $i = 1, 2, 3$, after receiving the pair $(W_{2,1}^{(1)}, W_{3,1}^{(2)})$.

With the second part, each user can obtain the subfiles of its desired file that have been cached by another user with the same demand. For example, the contents $W_{1,1} \oplus W_{1,2}$ and $W_{1,2} \oplus W_{1,3}$ help users U_1, U_2 and U_3 to obtain the subfiles of their request, W_1 , which have been cached by each other.

Finally, the last part of the delivery phase enables each user to decode the missing pieces of its requested file having been cached by another user with a different demand. For example, the delivered contents $W_{2,1}^{(2)} \oplus W_{2,2}^{(2)}, W_{2,2}^{(2)} \oplus W_{2,3}^{(2)}$, and $W_{1,4} \oplus W_{2,3}^{(2)}$ help users U_1, U_2 , and U_3 to obtain the piece $W_{1,4}^{(2)}$, and user U_4 can also decode the pieces $W_{2,1}^{(2)}, W_{2,2}^{(2)}$, and $W_{2,3}^{(2)}$. It can be verified that having received all the bits sent in three parts in the delivery phase, each user can obtain its desired file with a total delivery rate of $R_c = 2.1$. On the other hand, the MNC scheme achieves a delivery rate of $R_b = 2.12$ for the setting under consideration. \square

In the sequel, we present the cache placement and delivery phases of the proposed scheme in the general setting, analyze its delivery rate, and compare it with the state-of-the-art. We will observe that its superiority over the MNC scheme is not limited to the particular setting in the above example.

A. Placement Phase

We first generate K non-overlapping subfiles, each of size F/K bits, for each file W_i , $\forall i$, denoted by $W_{i,1}, \dots, W_{i,K}$. Similarly to [7] we use coded placement; that is, contents are cached in XOR-ed form in the placement phase. However, unlike in [7], instead of XORing subfiles of all the files in the database, we XOR subfiles in pairs. In particular, the following contents are cached by user U_k , for $k = 1, \dots, K$, in the placement phase:

$$Z_k = \bigcup_{i=1}^{N-1} (W_{i,k} \oplus W_{i+1,k}). \quad (12)$$

Since each subfile has a size of F/K bits, the limited memory of each cache is filled completely by the proposed placement scheme. In this way, each subfile of all the files is cached by exactly one user in the XOR-ed form. Hence, the whole of each file can be found in the caches of the users across the network (in coded form).

B. Delivery Phase

Note that, in the proposed caching scheme all the database is stored across the caches of the users. Therefore, in the delivery phase, the server first transmits the appropriate subfiles so that each user can recover all the subfiles stored in its cache in XOR-ed form. Then, the server transmits XOR of contents that are available at two different users, where each content is requested by the other user. This, equivalently, enables the two users to exchange their contents. By appropriately pairing subfiles, the server guarantees that each user receives the subfiles of its requested file that have been cached by every other user in the system.

Without loss of generality, by re-ordering the users, it is assumed that the first K_1 users, referred to as the group G_1 , request file W_1 , the next K_2 users, which form the group G_2 , demand W_2 , and so on so forth. For notational convenience, we define

$$S_i \triangleq \sum_{l=1}^i K_l, \quad (13)$$

where we set $S_0 \triangleq 0$. Thus, the general-case user demands can be expressed as follows:

$$d_k = i, \quad S_{i-1} + 1 \leq k \leq S_i, \quad \text{for } i = 1, \dots, N. \quad (14)$$

It is illustrated in Appendix A that the delivery rate of the proposed coded caching scheme does not depend on K_i , $\forall i$, i.e., the proposed scheme is not affected by the popularity of the files, as long as $K_i > 0$. Therefore, when $N < K$, the worst-case user demands for the proposed scheme happens when each file is requested by at least one user, i.e., $K_i \geq 1$, for $i = 1, \dots, N$.

The proposed delivery phase is divided into three distinct parts, and the contents delivered in part i is denoted by X_i , for $i = 1, 2, 3$. Hence, $X = (X_1, X_2, X_3)$ is transmitted over the shared link in the delivery phase. The delivery phase algorithm is presented for the worst-case user demands when $N < K$, i.e., when there is at least one user requesting each file. The proposed delivery phase algorithm is then extended to all values of N and K for a generic user demand combination assumption by introducing a new variable N' as the total number of files requested by the users.

To symmetrize the contents transmitted in the delivery phase, this phase is performed by further partitioning each subfile; that is, each subfile $W_{i,j}$, $\forall i, j$, is divided into $(N - 1)$ distinct pieces $W_{i,j}^{(1)}, \dots, W_{i,j}^{(N-1)}$, each of length $F/(K(N - 1))$ bits. Considering these smaller pieces, the content placed in the cache of user U_k , for $k = 1, \dots, K$, can

be re-written as follows:

$$Z_k = \bigcup_{l=1}^{N-1} \bigcup_{i=1}^{N-1} \left(W_{i,k}^{(l)} \oplus W_{i+1,k}^{(l)} \right). \quad (15)$$

Algorithm 1 Part 1 of the delivery phase

```

1: procedure PER-USER CODING
2:   for  $i = 1, \dots, N$  do
3:     for  $k = S_{i-1} + 1, \dots, S_i$  do
4:       for  $j = 1, \dots, N$  and  $j \neq i$  do
5:          $m_{j,k} = \begin{cases} j, & j < i \\ j - 1, & j > i \end{cases}$ 
6:          $X_1 \leftarrow (X_1, W_{j,k}^{(m_{j,k})})$ 
7:       end for
8:     end for
9:   end for
10: end procedure

```

The first part of the delivery phase is stated in Algorithm 1. The main purpose of this part is to enable each user U_k , $\forall k$, to retrieve all the subfiles $W_{j,k}$, for $j = 1, \dots, N$, that have been cached in the cache of user U_k in XOR-ed form. We remark that according to the cache placement in (15), for each $l \in [N-1]$, by delivering only one of the pieces $W_{1,k}^{(l)}, \dots, W_{N,k}^{(l)}$, user U_k can recover all the pieces $W_{1,k}^{(l)}, \dots, W_{N,k}^{(l)}$, $\forall k$. Hence, each user requires a total of $(N-1)$ distinct pieces to recover all the subfiles placed in its cache in XOR-ed form. To perform an efficient and symmetric delivery phase, $(N-1)$ distinct pieces, which are in the cache of user U_k in XOR-ed form, corresponding to $(N-1)$ different subfiles of the files that are not requested by user U_k are delivered to that user. For example, for user U_1 requesting file W_1 , the pieces $(W_{2,1}^{(1)}, W_{3,1}^{(2)}, \dots, W_{N,1}^{(N-1)})$ are delivered by Algorithm 1. Accordingly, for user U_k that has requested file W_i , the pieces $(W_{1,k}^{(1)}, \dots, W_{i-1,k}^{(i-1)}, W_{i+1,k}^{(i)}, \dots, W_{N,k}^{(N-1)})$ are delivered over the shared link. Thus, each user U_k can recover all subfiles $W_{j,k}$, $\forall j \in [N]$, stored in its cache in XOR-ed form. Note also that, Algorithm 1 delivers at most one piece of each subfile over the shared link. In Algorithm 1, we denote the index of the piece of subfile $W_{j,k}$, that is delivered in part 1 of the delivery phase by $m_{j,k}$. We will later refer to these indexes in explaining the other parts of the delivery phase. Note that, the pieces $(W_{1,k}^{(1)}, \dots, W_{i-1,k}^{(i-1)}, W_{i+1,k}^{(i)}, \dots, W_{N,k}^{(N-1)})$ are targeted for user U_k in group G_i demanding file W_i , for $i = 1, \dots, N$, and $k = S_{i-1} + 1, \dots, S_i$. Accordingly, for $j \in [N] \setminus \{i\}$, we have

$$m_{j,k} = \begin{cases} j, & j < i, \\ j - 1, & j > i, \end{cases} \quad (16)$$

which results in $m_{j,k} \leq N-1$.

For example, in Example 1 given above, we have $(m_{1,4}, m_{1,5}) = (1, 1)$, $(m_{2,1}, m_{2,2}, m_{2,3}, m_{2,5}) = (1, 1, 1, 2)$, and $(m_{3,1}, m_{3,2}, m_{3,3}, m_{3,4}) = (2, 2, 2, 2)$.

Algorithm 2 presents the second part of the proposed delivery phase, which allows each user to obtain its missing

subfiles that have been cached by the other users in the same group. Note that, having received part 1 of the delivery phase, user U_j in group G_i , for $i = 1, \dots, N$, and $j = S_{i-1} + 1, \dots, S_i$, can recover subfile $W_{i,j}$. Algorithm 2 delivers $\bigcup_{k=S_{i-1}+1}^{S_i-1} (W_{i,k} \oplus W_{i,k+1})$, with which user U_j can recover all the subfiles $W_{i,S_{i-1}+1}, \dots, W_{i,S_i}$, i.e., the subfiles of file W_i placed in the caches of users in group G_i .

Algorithm 2 Part 2 of the delivery phase

```

1: procedure INTER-GROUP CODING
2:   for  $i = 1, \dots, N$  do
3:      $X_2 \leftarrow \left( X_2, \bigcup_{k=S_{i-1}+1}^{S_i-1} (W_{i,k} \oplus W_{i,k+1}) \right)$ 
4:   end for
5: end procedure

```

The last part of the proposed delivery phase is presented in Algorithm 3, with which each user can receive the missing pieces of its desired file that have been placed in the cache of users in other groups. We deliver these pieces by exchanging data between the users in different groups. Observe that, for each user in group G_i , for $i = 1, \dots, N$, one piece of the subfile of its requested file W_i which is available to the users in G_j , for $j = 1, \dots, N$, $j \neq i$, was delivered in the first part of the delivery phase. Therefore, there are $(N-2)$ missing pieces of a file requested by a user, which have been placed in the cache of a user in a different group. For example, by delivering the pieces $(W_{2,1}^{(1)}, W_{3,1}^{(2)}, \dots, W_{N,1}^{(N-1)})$ to user U_1 demanding file W_1 in part 1 of the delivery phase, each user in group G_j with demand W_j can obtain the piece $W_{j,1}^{(j-1)}$, for $j = 2, \dots, N$. Therefore, there are $(N-2)$ missing pieces of the files requested by the users in groups G_2, \dots, G_N , that are available in the cache of user U_1 . Consider exchanging data between each user U_p in group G_i (demanding file W_i) and each user U_q in group G_j (demanding file W_j), for $i = 1, \dots, N-1$ and $j = i+1, \dots, N$, where $p = S_{i-1} + 1, \dots, S_i$ and $q = S_{j-1} + 1, \dots, S_j$. The subfile cached by user U_p (U_q) requested by user U_q (U_p) is $W_{j,p}$ ($W_{i,q}$). According to (16), the index of the piece of subfile $W_{j,p}$ ($W_{i,q}$) delivered in the first part of the delivery phase is equal to m_{j,S_i} (m_{i,S_j}), $\forall p \in \{S_{i-1} + 1, \dots, S_i\}$ and $\forall q \in \{S_{j-1} + 1, \dots, S_j\}$. Hence, the indexes of the missing pieces of each user in group G_i (G_j) available in the cache of a user in group G_j (G_i) are $[N-1] \setminus \{m_{i,S_j}\}$ ($[N-1] \setminus \{m_{j,S_i}\}$). Let $\pi_1^{i,j}(\cdot)$ and $\pi_2^{i,j}(\cdot)$ be arbitrary permutations on sets $[N-1] \setminus \{m_{i,S_j}\}$ and $[N-1] \setminus \{m_{j,S_i}\}$, respectively, for $i = 1, \dots, N-1$ and $j = i+1, \dots, N$. For $m_1 = \pi_1^{i,j}(l)$ and $m_2 = \pi_2^{i,j}(l)$, $\forall l \in [N-2]$, after receiving the corresponding contents delivered by Algorithm 3, all the users in G_i can recover the pieces $W_{i,S_{j-1}+1}^{(m_1)}, \dots, W_{i,S_j}^{(m_1)}$, and all the users in G_j can recover the pieces $W_{j,S_{i-1}+1}^{(m_2)}, \dots, W_{j,S_i}^{(m_2)}$.

Having received all three parts of the delivery phase, each user U_k , $\forall k$, can recover all the pieces of its desired file W_{d_k} that have been placed in any of the caches in the system. Together with the proposed placement phase, which guarantees

Algorithm 3 Part 3 of the delivery phase

```

1: procedure INTRA-GROUP CODING
2:   for  $i = 1, \dots, N - 1$  do
3:     for  $j = i + 1, \dots, N$  do
4:       for  $l = 1, \dots, N - 2$  do
5:          $m_1 = \pi_1^{i,j}(l)$ 
6:          $m_2 = \pi_2^{i,j}(l)$ 
7:          $X_3 \leftarrow \left( X_3, \bigcup_{k=S_{j-1}+1}^{S_j-1} \left( W_{i,k}^{(m_1)} \oplus W_{i,k+1}^{(m_1)} \right), \right.$ 
            $\left. \bigcup_{k=S_{i-1}+1}^{S_i-1} \left( W_{j,k}^{(m_2)} \oplus W_{j,k+1}^{(m_2)} \right), W_{i,S_j}^{(m_1)} \oplus W_{j,S_i}^{(m_2)} \right)$ 
8:       end for
9:     end for
10:  end for
11: end procedure

```

that all the subfiles of each file is available in one of the caches across the network, we can conclude that the demand of each user is satisfied by the proposed caching algorithm. It is to be noted that when $N = 2$, the proposed scheme is equivalent to the one proposed in [7], so we consider $N \geq 3$ throughout this paper.

C. Delivery Rate Analysis

The delivery rate of the proposed coded caching scheme is provided in the following theorem, whose detailed proof can be found in Appendix A.

Theorem 1. *In centralized caching with N files, each of length F bits, and K users, each equipped with a cache of capacity MF bits, if $N < K$ and $M = (N - 1)/K$, the following worst-case delivery rate is achievable:*

$$R_c \left(\frac{N-1}{K} \right) = N \left(1 - \frac{N}{2K} \right). \quad (17)$$

Remark 1. *To perform the proposed delivery phase for the general case, without loss of generality, the user demand combination is assumed as in (14), such that $K_i \geq 1$, for $i \leq N'$; and $K_i = 0$, otherwise, for some $N' \leq N$, that is a total of N' files are requested by the users in the system. In this case, each subfile is divided into $(N' - 1)$ distinct equal-length pieces, and in the delivery phase algorithm, the value N is substituted by N' . Hence, according to the delivery rate analysis provided in Appendix A, all the users' demands can be satisfied by delivering a total number of $R_c = N'(1 - N'/(2K))$ file(s). Since R_c is an increasing function of N' , we can conclude that, for $N < K$, the worst-case user demands happens when all the N files in the database are requested by at least one user, i.e., $K_i \geq 1$, for $i = 1, \dots, N$.*

Remark 2. *Based on Remark 1, when $N \geq K$, the worst-case user demands corresponds to the case when $N' = K$, i.e., all the users request distinct files in the database. Hence, the proposed scheme achieves a delivery rate of $K/2$ when $M = (N - 1)/K$, which is equal to the delivery rate of the*

state-of-the-art for the same cache capacity when $N = K$ and $N = K + 1$. However, the scheme proposed in [5] achieves a delivery rate smaller than $K/2$ for $M = (N - 1)/K$, when $N \geq K + 2$.

Remark 3. *It is possible to show that the proposed scheme improves upon the MNC scheme for $M = (N - 1)/K$. It is proven in Appendix B that the delivery rate achieved by memory-sharing between the scheme presented in [7] for $M = 1/K$ and the scheme proposed here for $M = (N - 1)/K$ has a smaller slope compared to the delivery rate of the MNC scheme at $M = 1/K$, when $K > N \geq 3$. Since the MNC scheme is achieved through memory-sharing between the schemes proposed in [7] and [5] for $M = 1/K < (N - 1)/K$ and $M = t^*N/K > (N - 1)/K$, respectively, where $t^* \geq 1$ is determined by (8), it is concluded that $R_c((N - 1)/K) \leq R_b((N - 1)/K)$, when $K > N \geq 3$.*

Remark 4. *We remark that the gain of the proposed scheme compared to the MNC scheme is due to the better use of the available cache capacities of the users when $MK = N - 1$, and partitioning each subfile into a number of distinct pieces, which allows performing an efficient and symmetric delivery phase. The proposed coded cache placement allows each user to retrieve all the subfiles in its cache at a relatively small cost. Furthermore, each user receives the bits of all the requested files rather than receiving bits only of its requested file, creating symmetry across users, which is helpful for the later steps of the delivery phase, and leads to a reduction in the delivery rate.*

Since $1/K < (N - 1)/K < t^*N/K$, the improvement of the proposed coded caching scheme for $M = (N - 1)/K$ can be extended to the range of cache capacities $M \in (1/K, t^*N/K)$ through memory-sharing with the caching schemes proposed in [7] and [5] for $M = 1/K$ and $M = t^*N/K$, respectively. In the following corollary, the improved delivery rate-cache capacity trade-off for $1/K \leq M \leq t^*N/K$ is presented.

Corollary 1. *The delivery rate-cache capacity trade-off*

$$R_c(M) = \begin{cases} -N \left(\frac{1}{2}M - 1 + \frac{1}{2K} \right), & \frac{1}{K} \leq M \leq \frac{N-1}{K}, \\ \frac{1}{(t^*-1)N+1} \left(\frac{K(K-t^*)}{t^*+1} - N \left(K - \frac{N}{2} \right) \right) \\ \quad \left(M - \frac{N-1}{K} \right) + N - \frac{N^2}{2K}, & \frac{N-1}{K} \leq M \leq \frac{t^*N}{K} \end{cases} \quad (18)$$

is achievable in a centralized caching system with a database of N files, and $K > N$ users, each equipped with a cache of normalized capacity M .

In the following theorem, which is proved in Appendix C, we show that the delivery rate-cache capacity trade-off $R_c(M)$ is within a constant multiplicative factor of 2 of the optimal delivery rate $R^*(M)$ for $1/K \leq M \leq (N - 1)/K$, when $K > N \geq 3$.

Theorem 2. *For a caching system with N files, and K users, satisfying $K > N \geq 3$, and a cache capacity of*

$M \in [1/K, (N-1)/K]$, we have

$$\frac{R_c(M)}{R^*(M)} \leq 2, \quad (19)$$

where $R_c(M)$ is the delivery rate achieved by the proposed coded caching scheme.

For all values of N and K , and all cache capacities $0 \leq M \leq N$, the delivery rate of the centralized coded caching scheme proposed in [5] is shown to be within a constant factor of 12 and 8 of the optimal delivery rate by utilizing the lower bounds derived in [5] and [10], respectively. Therefore, the proposed caching scheme reduces the best multiplicative gap in the literature by a factor of 4 for cache capacities satisfying $1/K \leq M \leq (N-1)/K$, when $K > N \geq 3$, achieved through memory-sharing between the proposed scheme for $M = (N-1)/K$ and the scheme of [7] for $M = 1/K$. Note that, the centralized coded caching scheme studied in [7] is optimal for cache capacities $M \leq 1/K$, when $N \leq K$.

D. Numerical Comparison with the State-of-the-Art

In this section, the delivery rate of the proposed scheme is compared numerically with the state-of-the-art results. In Fig. 3, the delivery rate of the proposed scheme for $M = (N-1)/K$, i.e., $R_c((N-1)/K)$ given by (17), is compared with that of the state-of-the-art for the same cache capacity, i.e., $R_b((N-1)/K)$ given by (9), as a function of $K \in [101, 1000]$ when $N = 100$. It can be seen that for the whole range of K values under consideration, the proposed coded caching scheme outperforms the MNC scheme, and the improvement is more noticeable for relatively moderate values of K . We also include in the figure two lower bounds on the delivery rate, the bound derived in [10, Theorem 1] and the cut-set based lower bound [5, Theorem 2]. Despite the improvement, there is still a large gap to the lower bounds. We believe that this gap is largely due to the looseness of the lower bound, but further improvements on the achievable delivery rate may also be possible. We are currently working on reducing this gap in both directions.

In Fig. 4, we compare the delivery rate-cache capacity trade-off achieved by the proposed coded caching strategy, $R_c(M)$, with the trade-off achieved by the MNC scheme, $R_b(M)$, when $N = 60$ and $K = 130$. In the figure, we focus on the cache capacity values $1/K \leq M \leq t^*N/K$ for which the proposed scheme outperforms the state-of-the-art. Note that, for this setting, based on (8), we have $t^* = 3$. Observe that the proposed scheme requires less data to be transmitted by the server over the shared link in the delivery phase for all cache capacity values satisfying $1/K < M < 3N/K$. We also include in the figure the two lower bounds on the delivery rate derived in [10, Theorem 1] and [5, Theorem 2].

IV. CONCLUSIONS

We have proposed a novel centralized coded caching scheme for a content delivery network consisting of a server delivering N popular files, each of length F bits, to K users, each with a cache of capacity MF bits. The proposed coded caching strategy, which is valid for all values of N and K , utilizes

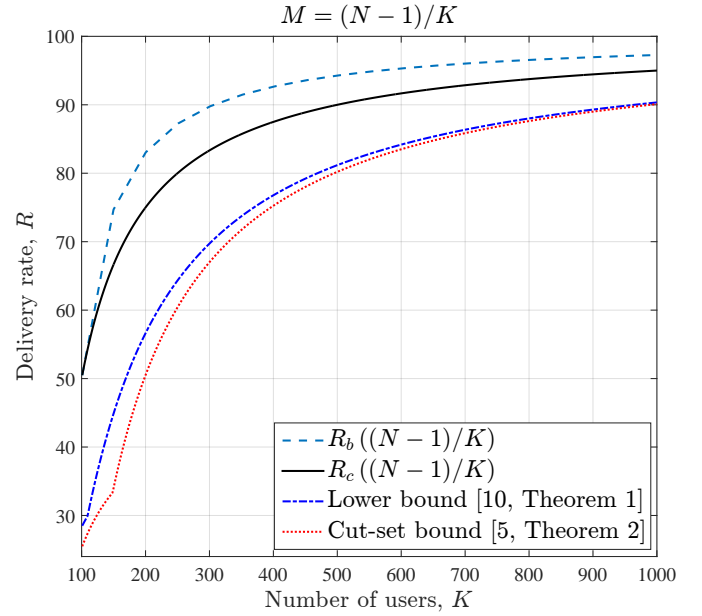


Fig. 3. Delivery rate for a caching system with $N = 100$ files as a function of the number of users, $K \in [101, 1000]$ when $M = (N-1)/K$.

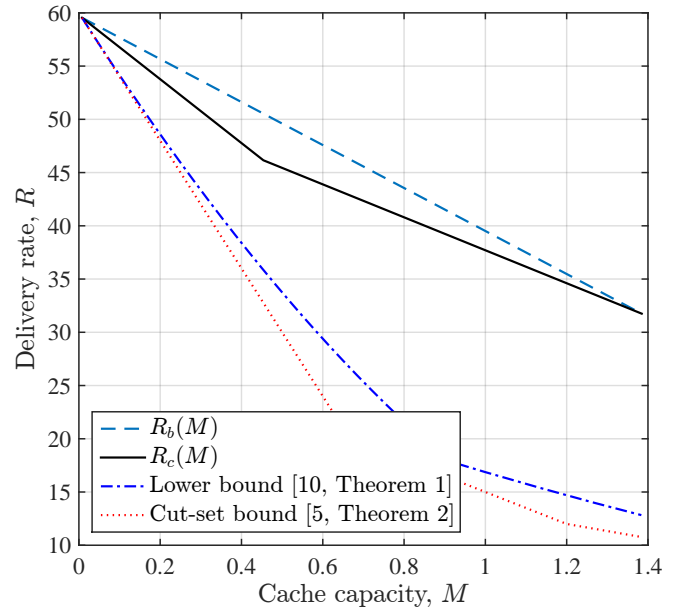


Fig. 4. Delivery rate-cache capacity trade-off for a caching system with $N = 60$ files and $K = 130$ users when $1/K \leq M \leq t^*N/K$, for which, according to (8), $t^* = 3$.

coded delivery, and creates symmetry between the delivered portions of different files in the delivery phase by partitioning files into different pieces. The delivery phase exploits both coded and uncoded transmission of various pieces of contents, carefully created to retain the symmetry across users and files. The delivery rate achieved by the proposed scheme for a cache capacity of $M = (N-1)/K$ is given by $R_c = N(1 - N/(2K))$ for $N < K$, which is shown to be lower than the state-of-the-art obtained by memory-sharing between the coded caching schemes proposed in [5] and [7]. We have then extended the improvement to a larger

range of cache capacities through memory-sharing between the proposed scheme and the best achievable scheme in the literature to obtain an order-optimal achievable delivery rate, which is shown to be within a constant multiplicative factor of 2 of the theoretically optimal delivery rate for cache capacities satisfying $1/K \leq M \leq (N-1)/K$, when $K > N \geq 3$.

We note that the caching scheme proposed in this paper places coded contents in the user cache memories, similarly to the scheme proposed in [7], and unlike the uncoded cache placement used in [5] and all the other follow-up works in the literature. We observe that, if the number of users in the system is more than the number of files, coded cache placement outperforms uncoded cache placement when the cache capacities are limited. We also note that, in the coded delivery formulation considered here, the total number of bits that need to be delivered over the shared link scales with F , the size of each content, which is assumed to be very large in this work. Therefore, the obtained gain in the delivery phase can be significant in terms of the number of delivered bits.

APPENDIX A PROOF OF THEOREM 1

To find the delivery rate of the proposed scheme, the delivery rate for each part of the delivery phase is calculated separately.

Having received the bits sent in the first part of the delivery phase presented in Algorithm 1, we would like each user to recover all the subfiles in its cache that have been cached in the XOR-ed form during the placement phase. However, to achieve this, we transmit pieces of the files that are not requested by that user. For example, for user U_k in group G_i with demand W_i , $i = 1, \dots, N$ and $k = S_{i-1} + 1, \dots, S_i$, we deliver $(N-1)$ different pieces corresponding to $(N-1)$ different files (except file W_i) to retrieve all the subfiles $W_{l,k}$, for $l = 1, \dots, N$. Since there are K users, a total of $K(N-1)$ different pieces, each of length $\frac{F}{K(N-1)}$ bits, are sent over the shared link in the first part of the delivery phase. As a result, the delivery rate of part 1 of the delivery phase is $R_{c_1} = 1$.

In part 2 of the proposed delivery phase provided in Algorithm 2, for the users in each group G_i , $(K_i - 1)$ XOR-ed contents $\bigcup_{k=S_{i-1}+1}^{S_i-1} (W_{i,k} \oplus W_{i,k+1})$ are transmitted over the shared link, enabling all the users in group G_i to recover the subfiles $W_{i,S_{i-1}+1}, \dots, W_{i,S_i}$. Hence, a total of $\sum_{i=1}^N (K_i - 1)$ XOR-ed contents, each of size F/K bits, are delivered over the shared link, which results in a delivery rate of

$$R_{c_2} = \frac{1}{K} \sum_{i=1}^N (K_i - 1) = 1 - \frac{N}{K} \quad (20)$$

for the second part of the delivery phase.

Finally, Algorithm 3 corresponds to the last part of the proposed delivery scheme, which enables file exchanges between the users in groups G_i and G_j , for $i = 1, \dots, N-1$ and $j = i+1, \dots, N$. There are $(N-2)$ missing pieces of the file requested by users in group G_i (G_j) that are located in the cache of each of the users in group G_j (G_i) with

indexes $l_1 \in [N-1] \setminus \{m_{i,S_j}\}$ ($l_2 \in [N-1] \setminus \{m_{j,S_i}\}$). Note that, we have $(N-2)$ missing pieces rather than $(N-1)$ as one piece was delivered in part 1 of the delivery scheme. For the piece with index l_1 and the piece with index l_2 , the server delivers $\bigcup_{k=S_{j-1}+1}^{S_j-1} (W_{i,k}^{(l_1)} \oplus W_{i,k+1}^{(l_1)})$,

$\bigcup_{k=S_{i-1}+1}^{S_i-1} (W_{j,k}^{(l_2)} \oplus W_{j,k+1}^{(l_2)})$, and $W_{i,S_j}^{(l_1)} \oplus W_{j,S_i}^{(l_2)}$, which enables all the users in group G_i to recover the pieces $W_{i,S_{j-1}+1}^{(l_1)}, \dots, W_{i,S_j}^{(l_1)}$, and also all the users in group G_j to recover the pieces $W_{j,S_{i-1}+1}^{(l_2)}, \dots, W_{j,S_i}^{(l_2)}$, by delivering a total of $(K_i + K_j - 1)$ XOR-ed contents, each of size $\frac{F}{K(N-1)}$ bits. As a result, the delivery rate of the third part is given by

$$R_{c_3} = \frac{(N-2)}{K(N-1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N (K_i + K_j - 1) = (N-2) \left(1 - \frac{N}{2K}\right). \quad (21)$$

By adding up the delivery rate of the three parts, the following delivery rate is achieved:

$$R_c \left(\frac{N-1}{K} \right) = R_{c_1} + R_{c_2} + R_{c_3} = N \left(1 - \frac{N}{2K}\right), \quad (22)$$

which completes the proof of Theorem 1.

APPENDIX B PROOF OF $R_c((N-1)/K) \leq R_b((N-1)/K)$

Since both $R_c(M)$, given by (18), and $R_b(M)$, given by (9), starts from $M = 1/K$ with the same rate of $N - N/K$ (utilizing the scheme proposed in [7]), to show that $R_c((N-1)/K) \leq R_b((N-1)/K)$ it suffices to prove that the slope of $R_c(M)$ at $M = 1/K$ is not larger than that of $R_b(M)$ at the same point. Observe that the slopes of $R_c(M)$ and $R_b(M)$ at $M = 1/K$ are $-N/2$ and α_{t^*} , determined according to (18) and (9), respectively. Due to the difficulty of characterizing α_{t^*} explicitly, we instead show that

$$\alpha_t \geq -\frac{N}{2}, \quad \forall t \in [K], \quad (23)$$

which concludes that $\alpha_{t^*} \geq -N/2$. We first define a function $g : ((\mathcal{Z}^+, \mathcal{Z}^+, [K]) \rightarrow \mathcal{R})$, where \mathcal{R} denotes the set of real numbers, as follows:

$$g(N, K, t) \triangleq \alpha_t + \frac{N}{2}, \quad (24)$$

and the goal is to illustrate that $g(N, K, t) \geq 0$, which is equivalent to

$$K(K-t) + N(t+1) \left(\frac{tN}{2} - K + \frac{1}{2} \right) \geq 0. \quad (25)$$

Inequality (25) can be re-written as follows:

$$\left(K - \frac{1}{2}(t + N(t+1)) \right)^2 - \frac{1}{4}(t + N(t+1))^2 + \frac{1}{2}N(t+1)(tN+1) \geq 0, \quad (26)$$

and after some algebraic manipulations, we have $g(N, K, t) \geq 0$, if and only if (iff)

$$h(N, K, t) \triangleq \left(K - \frac{1}{2}(t + N(t+1))\right)^2 + \frac{1}{4}(t^2 - 1)N(N-2) - \frac{1}{4}t^2 \geq 0. \quad (27)$$

Observe that for $t = 1$, we have

$$h(N, K, 1) = \left(K - N - \frac{1}{2}\right)^2 - \frac{1}{4} \geq 0. \quad (28)$$

Furthermore, for $t \geq 2$ and $N \geq 3$, we have

$$h(N, K, t) \geq \frac{1}{4}(t^2 - 1)N(N-2) - \frac{1}{4}t^2 \geq 0, \quad (29)$$

Consequently, we have

$$g(N, K, t) \geq 0, \quad \text{for } N \geq 3, \quad (30)$$

which results in $R_c((N-1)/K) \leq R_b((N-1)/K)$ for $N \geq 3$.

APPENDIX C PROOF OF THEOREM 2

We consider two distinct cases, when N is an even number, and when it is an odd number, and prove that for both cases, the multiplicative factor between the proposed achievable delivery rate $R_c(M)$ and $R^*(M)$ is at most 2 when $K > N \geq 3$, and $1/K \leq M \leq (N-1)/K$. According to the lower bound on the delivery rate derived in [10, Theorem 1], we have

$$R^*(M) \geq R_{LB}(M) \triangleq \max_{\substack{s \in \{1, \dots, K\}, \\ l \in \{1, \dots, \lceil \frac{N}{s} \rceil\}}} \frac{1}{l} \left\{ N - sM - \frac{\mu(N-ls)^+}{s+\mu} - (N-Kl)^+ \right\}, \quad (31)$$

where $\mu \triangleq \min\{\lceil (N-ls)/l \rceil, K-s\}$, and $(x)^+ \triangleq \max\{x, 0\}$.

First, assume that N is an even number. By setting $s = N/2$ and $l = 1$ in (31), for $K > N$, we have

$$R^*(M) \geq N - \frac{N}{2}M - \frac{\mu N}{N+2\mu}, \quad (32)$$

where

$$\mu = \min\left\{\left\lceil \frac{N}{2} \right\rceil, K - \frac{N}{2}\right\} = \frac{N}{2}, \quad (33)$$

which follows since $K > N$ and N is even. By substituting μ from (33) to (32), we have

$$R^*(M) \geq \frac{N}{4}(3-2M). \quad (34)$$

According to (18), for $1/K \leq M \leq (N-1)/K$, we have

$$\frac{R_c(M)}{R^*(M)} \leq \frac{2(2-M-1/K)}{3-2M}. \quad (35)$$

Note that, the expression on the right hand side of inequality (35) is an increasing function of M for $K \geq 2$. Setting the

cache capacity to its maximum value, $(N-1)/K$, we have

$$\frac{R_c(M)}{R^*(M)} \leq \frac{2(2K-N)}{3K-2N+2} \leq \frac{2(2K-N)}{3K-2N}. \quad (36)$$

The last expression above is a decreasing function of K , so by letting $K = N+1$, we have

$$\frac{R_c(M)}{R^*(M)} \leq \frac{2(N+2)}{N+3} \leq 2. \quad (37)$$

Now, consider N is an odd number. For $K > N$, we can set $s = (N-1)/2$ and $l = 1$ in the lower bound of (31) to find

$$R^*(M) \geq N - \frac{N-1}{2}M - \frac{(N+1)\mu}{N-1+2\mu}, \quad (38)$$

where

$$\mu = \min\left\{\left\lceil N - \frac{N-1}{2} \right\rceil, K - \frac{N-1}{2}\right\} = \frac{N+1}{2}. \quad (39)$$

Thus, we have

$$R^*(M) \geq N - \frac{N-1}{2}M - \frac{(N+1)^2}{4N}. \quad (40)$$

For $1/K \leq M \leq (N-1)/K$, we can obtain

$$\frac{R_c(M)}{R^*(M)} \leq \frac{2-M-1/K}{2-(1-\frac{1}{N})M-\frac{1}{2}(1+\frac{1}{N})^2}. \quad (41)$$

In the following, we show that the function $f : ([1/K, (N-1)/K]) \rightarrow \mathcal{R}$, defined as

$$f(M) \triangleq \frac{2-M-1/K}{2-(1-\frac{1}{N})M-\frac{1}{2}(1+\frac{1}{N})^2}, \quad (42)$$

is an increasing function of M for $K > N \geq 3$. We have

$$\frac{df}{dM} = \frac{(1-\frac{1}{N})\left(\frac{(N-1)K-2N}{2NK}\right)}{\left(2-(1-\frac{1}{N})M-\frac{1}{2}(1+\frac{1}{N})^2\right)^2} \geq 0, \quad (43)$$

where the last inequality in (43) holds for $K > N \geq 3$. Hence, we have

$$\begin{aligned} \frac{R_c(M)}{R^*(M)} &\leq f\left(\frac{N-1}{K}\right) \\ &= \frac{2-N/K}{2-(1-\frac{1}{N})\left(\frac{N-1}{K}\right)-\frac{1}{2}(1+\frac{1}{N})^2}. \end{aligned} \quad (44)$$

Now the goal is to prove that $f((N-1)/K) \leq 2$, for $K > N \geq 3$, when N is an odd number. After some simplification, it can be seen that $f((N-1)/K) \leq 2$ iff

$$p(N, K) \triangleq 1 - \frac{N-4}{K} - \frac{2}{KN} - \frac{1}{N^2} - \frac{2}{N} \geq 0. \quad (45)$$

Observe that, for $N \geq 4$, $p(N, K)$ is an increasing function of K . Thus, replacing $K = N+1$, we have

$$p(N, K) \geq \frac{N(3N-5)-1}{N^2(N+1)} \geq 0, \quad (46)$$

where the above inequality follows since $N \geq 4$. To complete the proof we need to show that $p(3, K) \geq 0$, for $K \geq 4$. We have

$$p(3, K) = \frac{2}{9} + \frac{1}{3K} \geq 0, \quad (47)$$

which completes the proof for odd N values. As a result, for $1/K \leq M \leq (N-1)/K$, the multiplicative gap between the delivery rate of the proposed scheme and the optimal delivery rate is at most 2 for all $K > N \geq 3$.

REFERENCES

- [1] L. W. Dowdy and D. V. Foster, "Comparative models of the file assignment problem," *ACM Comput. Surv.*, vol. 14, pp. 287–313, Jun. 1982.
- [2] K. C. Almeroth and M. H. Ammar, "The use of multicast delivery to provide a scalable and interactive video-on-demand service," *IEEE J. Sel. Areas Commun.*, vol. 14, no. 6, pp. 1110–1122, Aug. 1996.
- [3] S. Borst, V. Gupta, and A. Walid, "Distributed caching algorithms for content distribution networks," in *Proc. IEEE INFOCOM*, San Diego, CA, Mar. 2010, pp. 1–9.
- [4] A. C. Güngör and D. Gündüz, "Proactive wireless caching at mobile user devices for energy efficiency," in *Proc. IEEE Int'l Symp. on Wireless Comm. Systems (ISWCS)*, Brussels, Belgium, Aug. 2015.
- [5] M. A. Maddah-Ali and U. Niesen, "Fundamental limits of caching," *IEEE Trans. Inform. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [6] K. Wan, D. Tuninetti, and P. Piantanida, "On the optimality of uncoded cache placement," *arXiv: 1511.02256v1 [cs.IT]*, Nov. 2015.
- [7] Z. Chen, P. Fan, and K. B. Letaief, "Fundamental limits of caching: Improved bounds for small buffer users," *arXiv: 1407.1935v2 [cs.IT]*, Jul. 2014.
- [8] K. Wan, D. Tuninetti, and P. Piantanida, "On caching with more users than files," *arXiv: 1601.06383v2 [cs.IT]*, Jan. 2016.
- [9] M. Mohammadi Amiri and D. Gündüz, "Improved delivery rate-cache capacity trade-off for centralized coded caching," to be appeared in *Int'l Symp. on Inform. Theory and Its Applications (ISITA)*, Monterey, CA, Oct. 2016.
- [10] A. Sengupta, R. Tandon, and T. C. Clancy, "Improved approximation of storage-rate tradeoff for caching via new outer bounds," in *Proc. IEEE Int'l Symp. on Inform. Theory (ISIT)*, Hong Kong, Jun. 2015, pp. 1691–1695.
- [11] H. Ghasemi and A. Ramamoorthy, "Improved lower bounds for coded caching," in *Proc. IEEE Int'l Symp. on Inform. Theory (ISIT)*, Hong Kong, Jun. 2015, pp. 1696–1700.
- [12] C. Tian, "A note on the fundamental limits of coded caching," *arXiv: 1503.00010v1 [cs.IT]*, Feb. 2015.
- [13] M. A. Maddah-Ali and U. Niesen, "Decentralized caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Trans. Netw.*, vol. 23, no. 4, pp. 1029–1040, Apr. 2014.
- [14] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, and S. Diggavi, "Hierarchical coded caching," in *Proc. IEEE Int'l Symp. on Inform. Theory (ISIT)*, Honolulu, HI, Jun. 2014, pp. 2142–2146.
- [15] J. Zhang, X. Lin, C. C. Wang, and X. Wang, "Coded caching for files with distinct file sizes," in *Proc. IEEE Int'l Symp. on Inform. Theory (ISIT)*, Hong Kong, Jun. 2015, pp. 1686–1690.
- [16] U. Niesen and M. A. Maddah-Ali, "Coded caching with nonuniform demands," in *Proc. IEEE Conf. Comput. Commun. Workshops (INFOCOM WKSHPS)*, Toronto, ON, Apr. 2014, pp. 221–226.
- [17] M. Ji, A. M. Tulino, J. Llorca, and G. Caire, "Order-optimal rate of caching and coded multicasting with random demands," *arXiv: 1502.03124v1 [cs.IT]*, Feb. 2015.
- [18] S. Wang, W. Li, X. Tian, and H. Liu, "Coded caching with heterogeneous cache sizes," *arXiv:1504.01123v3 [cs.IT]*, Aug. 2015.
- [19] M. Mohammadi Amiri, Q. Yang, and D. Gündüz, "Decentralized coded caching with distinct cache capacities," *arXiv:1610.03792v1 [cs.IT]*, Oct. 2016.
- [20] K. Poularakis and L. Tassioulas, "On the complexity of optimal content placement in hierarchical caching networks," *IEEE Trans. Commun.*, vol. 64, no. 5, pp. 2092–2103, May 2016.
- [21] Q. Yang and D. Gündüz, "Centralized coded caching for heterogeneous lossy requests," in *Proc. IEEE Int'l Symp. on Inform. Theory (ISIT)*, Barcelona, Spain, Jul. 2016.
- [22] P. Hassanzadeh, E. Erkip, J. Llorca, and A. Tulino, "Distortion-memory tradeoffs in cache-aided wireless video delivery," *arXiv:1511.03932v1 [cs.IT]*, Nov. 2015.
- [23] R. Timo, S. B. Bidokhti, M. Wigger, and B. Geiger, "A rate-distortion approach to caching," in *Proc. International Zurich Seminar on Communications*, Zurich, Switzerland, Mar. 2016.
- [24] R. Pedarsani, M. A. Maddah-Ali, and U. Niesen, "Online coded caching," in *Proc. IEEE Int'l Conf. Commun. (ICC)*, Sydney, Australia, Jun. 2014, pp. 1878–1883.
- [25] R. Timo and M. Wigger, "Joint cache-channel coding over erasure broadcast channels," *arXiv:1505.01016v1 [cs.IT]*, May 2015.
- [26] S. Saeedi Bidokhti, M. Wigger, and R. Timo, "Noisy broadcast networks with receiver caching," *arXiv:1605.02317v1 [cs.IT]*, May 2016.
- [27] W. Huang, S. Wang, L. Ding, F. Yang, and W. Zhang, "The performance analysis of coded cache in wireless fading channel," *arXiv:1504.01452v1 [cs.IT]*, Apr. 2015.
- [28] M. Gregori, J. Gomez-Vilardebo, J. Matamoros, and D. Gündüz, "Wireless content caching for small cell and D2D networks," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 5, pp. 1222–1234, May 2016.
- [29] M. Ji, G. Caire, and A. F. Molisch, "Fundamental limits of caching in wireless D2D networks," *arXiv:1405.5336v1 [cs.IT]*, May 2014.
- [30] K. Shanmugam, N. Golrezaei, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless content delivery through distributed caching helpers," *IEEE Trans. Inform. Theory*, vol. 59, no. 12, pp. 8402–8413, Sep. 2013.
- [31] N. Golrezaei, A. F. Molisch, A. G. Dimakis, and G. Caire, "Femtocaching and device-to-device collaboration: A new architecture for wireless video distribution," *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, Apr. 2013.



processing.

Mohammad Mohammadi Amiri received the B.Sc. degree in Communication Systems from the Iran University of Science and Technology in 2011 and the M.Sc. degree in Communication Systems engineering from the University of Tehran in 2014 both with highest rank in classes. Currently, he is a research postgraduate at Imperial College London under the supervision of Dr. Gündüz. His research interests include information and coding theory, wireless communications, MIMO systems, cooperative networks, cognitive radio, and signal



processing.

Deniz Gündüz [S'03-M'08-SM'13] received the B.S. degree in electrical and electronics engineering from METU, Turkey in 2002, and the M.S. and Ph.D. degrees in electrical engineering from NYU Polytechnic School of Engineering in 2004 and 2007, respectively. After his PhD, he served as a postdoctoral research associate at Princeton University, and as a consulting assistant professor at Stanford University. He was a research associate at CTTC in Barcelona, Spain until September 2012, when he joined the Electrical and Electronic Engineering Department of Imperial College London, UK, as a Lecturer. Currently he is a Reader in the same department.

Dr. Gündüz is an Editor of the IEEE Transactions on Communications, and IEEE Transactions on Green Communications and Networking. He is the recipient of a Starting Grant of the European Research Council (ERC), the 2014 IEEE Communications Society Best Young Researcher Award for the Europe, Middle East, and Africa Region, Best Paper Award at the 2016 IEEE Wireless Communications and Networking Conference (WCNC), and the Best Student Paper Award at the 2007 IEEE International Symposium on Information Theory (ISIT). He served as the General Co-chair of the 2016 IEEE Information Theory Workshop, and the 2012 IEEE European School of Information Theory (ESIT). His research interests lie in the areas of communication theory and information theory with special emphasis on joint source-channel coding, multi-user networks, energy efficient communications and privacy in cyber-physical systems.